

WHEN DOES IT PAY TO BREAK THE MATCHES FOR ANALYSIS OF A MATCHED PAIRS DESIGN?

Henry S. Lynn, Charles E. McCulloch

Biometrics Unit, Cornell University, Ithaca, NY 14853

BU-1042-MB¹

January, 1990

SUMMARY

Two methods of analyses are compared to estimate the treatment effect of a comparative study where each treated individual is matched with a single control at the design stage. The usual matched pairs analysis accounts for the pairing directly in its model, whereas regression adjustment ignores the matching but instead models the pairing using a set of covariates. For a normal linear model, the estimated treatment effect from the matched pairs analysis (paired t-test) is more efficient. For a Bernoulli logistic model, matched pairs analysis performs better when the sample size was small, but is inferior to logistic regression for large sample sizes.

¹Technical Report BU-1042-MA in the Biometrics Unit Series.

1. INTRODUCTION

In many comparative studies, researchers may decide at the design stage to form pairs by matching exactly on some not so readily quantifiable variable; e.g., sibship or neighborhood of residence, but in the subsequent analysis it is not always clear what statistical analyses they should select. The usual approach is to employ a matched pairs analysis. This method will yield an unbiased or asymptotically unbiased estimate of the "treatment" effect, although the estimate may have a larger variance than, say, that obtained from a regression analysis which attempts to model the pairing using some set of covariates. The latter approach sacrifices the unbiasedness of the estimate in an attempt to gain higher precision. The purpose of this paper is then to determine, for both the normal linear model and Bernoulli logistic model, which of the above two analyses is preferable by quantitatively assessing the trade off between the loss in accuracy and increase in efficiency involved in modelling the pairing.

Prentice (1976), and Breslow and Day (1980), using examples of case-control studies, have shown that a logistic regression analysis which failed to account for "important" covariates will tend to bias the regression estimates, whereas an analysis that included "redundant" covariates will inflate the variances of the estimates. In particular, Prentice (1976) cites a matched pairs study relating post-menopausal estrogen exposure on endometrical cancer, and recommends an analysis without retaining the pairing since pairs members were not intrinsically similar in respects other than those indicated by the matching variables. However, their

findings were more qualitative, and do not indicate how "important" the covariates need to be; i.e., how adequate the covariates are in explaining the pairing before regression analysis becomes more effective than a matched pairs analysis. Extensive work on the problem of omitted covariates in general linear models has been done by Gail et al. (1984, 1988), in which the covariates are treated as random variables. Many authors have also examined the pros and cons of matching versus regression adjustment, although their comparisons are restricted to the case when the pairing has been fully modelled by the covariates. Others have approached the problem from an experimental viewpoint, comparing the efficiencies of unmatched versus matched designs. For an overview of some of the results, see Rubin (1973), McKinlay (1977), Kupper et al. (1981), and Greenland (1986).

In the following comparisons between matching and regression adjustment, it will be assumed that we have a comparative study involving a control group and a "treatment" group. There will be a total of m pairs and $n = 2m$ subjects, where each subject from the treatment group is paired with a single subject from the control group.

2. NORMAL LINEAR MODEL

Suppose that the true model is given by

$$Y = X\beta + Z\gamma + e, \quad (1)$$

where Y is arranged such that (y_i, y_{i+m}) , $1 \leq i \leq m$, are the responses of the i th pair, $\beta' = (\beta_1, \beta_2, \dots, \beta_s)$ and $\gamma' = (\gamma_1, \dots, \gamma_t)$ represent vectors of

parameters, $\mathbf{X} = \begin{bmatrix} \mathbf{1} & \mathbf{1} & | & \mathbf{X}_1 \\ \mathbf{1} & -\mathbf{1} & | & \mathbf{X}_2 \end{bmatrix}_{n \times s}$ and $\mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1 \\ \mathbf{Z}_2 \end{bmatrix}_{n \times t}$ are matrices of

arbitrary fixed constants, and $E(\mathbf{e}) = \mathbf{0}$ and $\text{Var}(\mathbf{e}) = \sigma_e^2 \mathbf{I}_n$. This model therefore assumes that the "intra-match" correlation can be explained by a finite set of covariates within a linear regression model format.

In the following, \mathbf{X} will represent the covariates used to model the pairing and \mathbf{Z} will represent the covariates omitted from the model. Moreover, \mathbf{X}_1 and \mathbf{Z}_1 are the covariates for the treatment group whereas \mathbf{X}_2 and \mathbf{Z}_2 are the covariates for the control group, and $\mathbf{X}_1 = \mathbf{X}_2$ and $\mathbf{Z}_1 = \mathbf{Z}_2$ because of exact matching. For the model stated in (1), the object of inference will be the treatment effect, defined as $2\beta_2$ for simplicity of analysis.

2.1) Matched Pairs Analysis.

The matched pairs analysis is the usual analysis for two paired samples, with the true model rewritten as

$$\mathbf{Y} = \begin{bmatrix} \mathbf{1} & \mathbf{1} \\ \mathbf{1} & -\mathbf{1} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} + \boldsymbol{\varepsilon} \quad , \quad \text{where} \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \end{bmatrix} \begin{bmatrix} \beta_3 \\ \vdots \\ \beta_s \end{bmatrix} + \begin{bmatrix} \mathbf{Z}_1 \\ \mathbf{Z}_2 \end{bmatrix} \gamma + \mathbf{e} \quad . \quad (2)$$

If we define the pairwise differences,

$$y_i - y_{m+i} = 2\beta_2 + (\varepsilon_i - \varepsilon_{m+i}) \quad \text{for } i = 1, 2, \dots, m \quad ,$$

the treatment effect is then estimated by $\hat{2\beta_2} \equiv \frac{1}{m} \sum_{i=1}^m (y_i - y_{m+i})$, which is

unbiased since $E(\hat{2\beta}_2) = 2\beta_2 + \frac{1}{m} \sum_{i=1}^m E(\epsilon_i - \epsilon_{m+i}) = 2\beta_2$. Also

$$\text{Var}(\hat{2\beta}_2) = \frac{1}{m} \text{Var}(\epsilon_i - \epsilon_{m+i}) = \frac{1}{m} \text{Var}(e_i - e_{m+i}) = \frac{2}{m} \sigma_e^2, \quad (3)$$

where σ_e^2 is estimated by $\hat{\sigma}_e^2$, which equals half the value of the estimated variance of the pairwise differences.

2.2) Regression Adjustment.

Suppose that in our study we have measured \mathbf{X} , a subset of the covariates that determine the pairing. We now decide to break the matches in the analysis and instead model the pairing by regressing \mathbf{Y} on \mathbf{X} . We adopt a working model under which

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\xi}, \quad (4)$$

where $E(\boldsymbol{\xi}) = \mathbf{0}$ and $\text{Var}(\boldsymbol{\xi}) = \sigma_\xi^2 \mathbf{I}_n$. Note that under the true model $\boldsymbol{\xi} = \mathbf{Z}\boldsymbol{\gamma} + \mathbf{e}$, such that $E(\boldsymbol{\xi}) = \mathbf{Z}\boldsymbol{\gamma}$ and $\text{Var}(\boldsymbol{\xi}) = \sigma_e^2 \mathbf{I}_n$. The usual least squares analysis then gives $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}$ with $\text{Var}(\hat{\boldsymbol{\beta}}) = \sigma_\xi^2 (\mathbf{X}'\mathbf{X})^{-1}$. Observe that the above is just an underfitted model of the true model in (1), with \mathbf{Z} being the omitted covariates. In general this would mean that $\hat{\boldsymbol{\beta}}$ is biased since

$$E(\hat{\boldsymbol{\beta}}) = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma}) = \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Z}\boldsymbol{\gamma}.$$

However in this case $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1} \begin{bmatrix} \mathbf{1}' & \mathbf{1}' \\ \mathbf{1}' & -\mathbf{1}' \\ \mathbf{X}'_1 & \mathbf{X}'_2 \end{bmatrix} \begin{bmatrix} \mathbf{Z}_1 \\ \mathbf{Z}_2 \end{bmatrix} \boldsymbol{\gamma}$

$$\begin{aligned}
 &= \beta + \begin{bmatrix} 2m & 0 & 1'(X_1 + X_2) \\ 0 & 2m & 0' \\ (X_1 + X_2)'1 & 0 & X_1'X_1 + X_2'X_2 \end{bmatrix}^{-1} \begin{bmatrix} 1' & 1' \\ 1' & -1' \\ X_1' & X_2' \end{bmatrix} \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} \gamma \\
 &= \beta + \begin{bmatrix} C_{11} & 0 & C_{12} \\ 0 & \frac{1}{2m} & 0' \\ C_{21} & 0 & C_{22} \end{bmatrix} \begin{bmatrix} 1' & 1' \\ 1' & -1' \\ X_1' & X_2' \end{bmatrix} \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix} \gamma,
 \end{aligned}$$

where $C = \begin{bmatrix} 2m & 1'(X_1 + X_2) \\ (X_1 + X_2)'1 & X_1'X_1 + X_2'X_2 \end{bmatrix}^{-1}$.

This implies that $E(\hat{\beta}_2) = \beta_2 + \frac{1}{2m}(1'Z_1 - 1'Z_2)\gamma = \beta_2$, and thus $2\hat{\beta}_2$ is again an unbiased estimator of the treatment effect because $Z_1 = Z_2$ due to exact matching.

Therefore to compare the effectiveness of the matched pairs analysis with regression adjustment, we consider the variance of $2\hat{\beta}_2$ since $2\hat{\beta}_2$ is unbiased for both models. Specifically, we will use the length of a 95% confidence interval of $2\hat{\beta}_2$ as our criterion of comparison. Our working model (4) implies using $(X'X)^{-1}$ above that $\text{Var}(\hat{\beta}_2) = \frac{\sigma_\xi^2}{2m}$, and hence

$$\text{Var}(2\hat{\beta}_2) = \frac{2\sigma_\xi^2}{m}, \quad (5)$$

where σ_ξ^2 is estimated by $\hat{\sigma}_\xi^2$, the residual mean square from the regression of Y on X . Thus, if we define λ as the ratio of the lengths of the 95% confidence intervals of $2\hat{\beta}_2$ from the regression adjustment over

the matched pairs analysis, then using (3) and (5) we have

$$\lambda = \frac{2 t_{2n-s} \sqrt{\frac{2}{m} \hat{\sigma}_{\xi}^2}}{2 t_{m-1} \sqrt{\frac{2}{m} \hat{\sigma}_e^2}}, \quad (6)$$

where t_v is the upper 0.025 percentage point of the t distribution with v degrees of freedom.

In order to determine the behavior of λ we will now evaluate $\hat{\sigma}_e^2$ and $\hat{\sigma}_{\xi}^2$ in terms of their unobservable theoretical values. According to the true model in (1), $\hat{\sigma}_e^2$ is the residual mean square from the regression of Y on X and Z . Comparing (1) with (4) we find that $\xi = Z\gamma + e$, and thus

$$(n-s) \hat{\sigma}_{\xi}^2 = SSR(Z|X) + (n-s-t) \hat{\sigma}_e^2, \quad (7)$$

where

$$SSR(Z|X) = Y'(PZ)((PZ)'PZ)^{-1}(PZ)'Y \quad (8)$$

is the sequential sum of squares of Z given that X has already been fitted, and P is defined to be equal to $I_n - X(X'X)^{-1}X'$.

Now observe that $PY = P(X\beta + \xi) = P\xi$, and since P is idempotent and symmetric this implies that

$$Y'PZ = (Y'P)PZ = (P\xi)'PZ. \quad (9)$$

Substituting (9) into (8), we obtain

$$SSR(Z|X) = (P\xi)'(PZ)((PZ)'PZ)^{-1}(PZ)'P\xi, \quad (10)$$

which means that $SSR(Z|X)$ is also the regression sum of squares of the regression of $P\xi$ on PZ . Applying this result, we can then rewrite (7) as

$$\hat{\sigma}_e^2 = \frac{n-s}{n-s-t} (1 - R^2) \hat{\sigma}_{\xi}^2, \quad (11)$$

where R^2 is the coefficient of determination of the regression of $P\xi$ on PZ .

Consequently, on substituting (11) into (6) we obtain

$$\lambda = \frac{t_{2n-s}}{t_{m-1} \sqrt{1 - R_a^2}},$$

where $1 - R_a^2 = \frac{n-s}{n-s-t} (1 - R^2)$.

We use the notation R_a^2 since the quantity defined is similar to an adjusted R^2 . In fact, if ξ were observable so that we could fit a model of the form $\xi = Z\gamma + e$, then it can be shown that

$$\hat{\sigma}_e^2 = (1 - R_a^2) \hat{\sigma}_\xi^2,$$

where R_a^2 is the adjusted R^2 for the regression of ξ on Z .

To illustrate the efficacy of the matched pairs analysis over regression adjustment, consider Table 1 which lists the values of λ for the simple case when $s = 3$ and $t = 1$; i.e. when there is one known covariate and one omitted covariate. Observe that the gain in precision for the regression analysis, due to its larger degrees of freedom, diminishes with increasing sample sizes and values of R_a^2 , and is never significantly more efficient than the matched pairs analysis for any reasonable sample size and value of R_a^2 . For example, if $m = 30$, then $\lambda = 1.03$ when $R_a^2 = 0.10$. Intuitively, this means that the measured covariate needs to explain more than 90% of the variation of the regression of $P\xi$ on PZ (or ξ on Z) before regression adjustment can be as effective as the matched pairs analysis. The behavior of λ follows similarly when there are more covariates; i.e., for larger values of s and t . In summary, the above results indicate that the matched pairs analysis is preferable for the normal linear model.

Table 1. Ratio of the length of a 95% confidence interval using a regression adjustment over a matched pairs analysis.

m	R_a^2					
	0.0	0.1	0.3	0.5	0.7	0.9
5	.85	.90	1.02	1.20	1.56	2.70
10	.93	.98	1.11	1.32	1.70	2.95
30	.98	1.03	1.17	1.39	1.79	3.10
50	.99	1.04	1.18	1.40	1.80	3.13
100	.99	1.05	1.19	1.41	1.81	3.14

3. BERNOULLI LOGISTIC MODEL

Suppose now that y_1, \dots, y_n are independent binary random variables with densities $h_i(y_i) = \pi_i^{y_i} (1 - \pi_i)^{1-y_i}$, and the true model is given by

$$\text{logit}(\pi_i) = \ln \frac{\pi_i}{1 - \pi_i} = \mathbf{x}_i' \boldsymbol{\beta} + \mathbf{z}_i' \boldsymbol{\gamma}, \quad (12)$$

$$\text{so } \pi_i = E(y_i) = \Pr\{y_i = 1\} = \frac{1}{1 + e^{-\mathbf{x}_i' \boldsymbol{\beta} - \mathbf{z}_i' \boldsymbol{\gamma}}},$$

$$\boldsymbol{\beta}' = (\beta_1, \dots, \beta_s), \quad \boldsymbol{\gamma}' = (\gamma_1, \dots, \gamma_t),$$

$$\mathbf{x}_i' \text{ is the } i\text{th row of } \mathbf{X} = \begin{bmatrix} \mathbf{1} & \mathbf{1} & | & \mathbf{X}_1 \\ \mathbf{1} & -\mathbf{1} & | & \mathbf{X}_2 \end{bmatrix}_{n \times s}, \text{ and}$$

$$\mathbf{z}_i' \text{ is the } i\text{th row of } \mathbf{Z} = \begin{bmatrix} \mathbf{Z}_1 \\ \mathbf{Z}_2 \end{bmatrix}_{n \times t}, \quad i = 1, 2, \dots, n.$$

As in Section 2, \mathbf{X}_1 and \mathbf{Z}_1 are the covariates for the treatment group whereas \mathbf{X}_2 and \mathbf{Z}_2 are the covariates for the control group. Furthermore, $\mathbf{X}_1 = \mathbf{X}_2$ and $\mathbf{Z}_1 = \mathbf{Z}_2$ because of exact matching, and $2\beta_2$ will again be used to measure the treatment effect.

3.1) Matched Pairs Analysis.

For a matched pairs design, the typical quantity of interest is the odds ratio, Ψ , or the log odds ratio $\tau = \ln \Psi$. According to model (12) we have $\Psi = e^{2\beta_2}$ and $\tau = 2\beta_2$. Hence, if we let $n_{10} \equiv$ number of pairs with a "1" for the treated subject and a "0" for the control, $n_{01} \equiv$ number of pairs with a "0" for the treated subject and a "1" for the control, and $N_d = n_{10} + n_{01} \equiv$ number of discordant pairs, then according to Breslow and Day (1981) the usual estimator for τ is

$$\hat{\tau}_{ML} = \ln \hat{\Psi} = \ln \left(\frac{n_{10}}{n_{01}} \right) = \ln \left(\frac{n_{10}}{N_d - n_{10}} \right), \quad (13)$$

which is the maximum likelihood estimator of the log odds ratio.

In order to evaluate the efficiency of $\hat{\tau}_{ML}$ and subsequently compare it with the estimator from the logistic regression analysis, we need to first obtain the expectation and variance of $\hat{\tau}_{ML}$. Applying a Taylor's series approximation to (13), it can be shown that

$$\text{Var}(\hat{\tau}_{ML} | N_d) \cong \frac{(\Psi + 1)^2}{N_d \Psi}, \text{ and} \quad (14)$$

$$E(\hat{\tau}_{ML} | N_d) \cong \tau + \frac{(\Psi + 1)(\Psi - 1)}{2N_d \Psi}. \quad (15)$$

It follows from (15) that

$$E(\hat{\tau}_{ML}) = E(E(\hat{\tau}_{ML} | N_d)) \cong \tau + \frac{(e^\tau + 1)(e^\tau - 1)}{2e^\tau} E\left(\frac{1}{N_d}\right), \quad (16)$$

and using (14) and the first term of $E(\hat{\tau}_{ML} | N_d)$ in (15) we then find that

$$\text{Var}(\hat{\tau}_{ML}) = E(\text{Var}(\hat{\tau}_{ML} | N_d)) + \text{Var}(E(\hat{\tau}_{ML} | N_d))$$

$$\cong E(\text{Var}(\hat{\tau}_{\text{ML}}|N_d)) \cong \frac{(e^{\tau} + 1)^2}{e^{\tau}} E\left(\frac{1}{N_d}\right). \quad (17)$$

Since the square of the bias is of order $E\left(\frac{1}{N_d}\right)^2$, it is negligible with respect to the variance and asymptotically we have

$$\text{MSE}(\hat{\tau}_{\text{ML}}) \cong \frac{(e^{\tau} + 1)^2}{e^{\tau}} E\left(\frac{1}{N_d}\right). \quad (18)$$

3.2) Logistic Regression.

Analogous to the situation described in Section 2.2, we assume here that, \mathbf{Z} , a subset of the covariates has been omitted, thus while the true model is (12), we are using a working model in which y_1, \dots, y_n are assumed to have densities $f_i(y_i, \boldsymbol{\beta}) = P_i^{y_i} (1 - P_i)^{1-y_i}$, where

$$P_i = \frac{1}{1 + e^{-\mathbf{x}_i' \boldsymbol{\beta}}} \quad \text{and} \quad \boldsymbol{\beta}' = (\beta_1, \dots, \beta_s).$$

Applying a result in Royall's (1986, p.222) paper, we then find that $\hat{\boldsymbol{\beta}}$, the maximum likelihood estimator of $\boldsymbol{\beta}$, is a consistent estimator of the root of the likelihood equation, $\boldsymbol{\beta}^o$. For large n , $\boldsymbol{\beta}^o$ is approximately the solution of $\sum_{i=1}^n (\pi_i - P_i) \mathbf{x}_i' = \mathbf{0}$, where the π_i 's are the true probabilities as defined by the model in (12).

$$\text{Furthermore, if we define } P_i^o = \frac{1}{1 + e^{-\mathbf{x}_i' \boldsymbol{\beta}^o}},$$

$$\mathbf{V} = \begin{bmatrix} \mathbf{V}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_2 \end{bmatrix}, \quad \text{where } \mathbf{V} = \text{diag}\{\pi_i (1 - \pi_i)\} \text{ for } i = 1, 2, \dots, n, \text{ and}$$

$$\mathbf{V}^o = \begin{bmatrix} \mathbf{V}_1^o & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_2^o \end{bmatrix}, \text{ where } \mathbf{V}^o = \text{diag}\{P_i^o(1 - P_i^o)\} \text{ for } i = 1, 2, \dots, n,$$

then (Royall, p.222, 1986)

$$\sqrt{n} (\hat{\beta} - \beta) \xrightarrow{L} \mathcal{N}(\beta^o - \beta, \Sigma), \quad (19)$$

where $\Sigma = n (\mathbf{X}' \mathbf{V}^o \mathbf{X})^{-1} (\mathbf{X}' \mathbf{V} \mathbf{X}) (\mathbf{X}' \mathbf{V}^o \mathbf{X})^{-1}$.

Therefore, on applying equation (19), we conclude that

$$\text{Bias}(2\hat{\beta}_2) = 2(\beta_2^o - \beta_2), \quad (20)$$

and

$$\text{Var}(2\hat{\beta}_2) = 4\Sigma_{22}, \quad (21)$$

where Σ_{22} is the second diagonal element of Σ .

3.3) Results and Discussion.

Recall that $\hat{\tau}_{\text{ML}}$ is the ML estimator of the treatment effect from the matched pairs analysis. Now let $\hat{\tau}_{\text{logit}}$ be the corresponding estimator obtained using logistic regression, then from (20) and (21) we find that asymptotically,

$$\text{MSE}(\hat{\tau}_{\text{logit}}) = 4\{\Sigma_{22} + (\beta_2^o - \beta_2)^2\}. \quad (22)$$

To compare the matched pairs analysis with logistic regression, we consider the simple case when there is only one known covariate and one omitted covariate; i.e. the true model is given by

$$\text{logit}(\pi_i) = \beta_1 + \beta_2 x_{2i} + \beta_3 x_{3i} + \gamma z_i, \quad 1 \leq i \leq n.$$

The comparison was investigated in terms of the asymptotic performances of $\hat{\tau}_{ML}$ and $\hat{\tau}_{logit}$ via equations (18) and (22). The quantity $E\left(\frac{1}{N_d}\right)$ in (18) was approximated either using $1/E(N_d)$ or through a normal approximation for $Pr(N_d = k)$, where $k=1,2,\dots,m$ (the case $k=0$ is omitted). Since the mean square errors are a complicated function of the parameters and the design, we calculated them for a 2^7 arrangement of the following factors:

$m \equiv$ number of pairs,

$x_3 \equiv$ known covariate vector,

$\beta_1 \equiv$ intercept,

$\beta_2 \equiv$ regression coefficient for considering the treatment effect,

$\beta_3 \equiv$ regression coefficient for the known covariate,

$\gamma \equiv$ regression coefficient for the omitted covariate, and

$\rho \equiv$ correlation between x_3 and z (the omitted covariate vector).

ρ was chosen as a measure of effectiveness of pairing since, intuitively, if the correlation is high, the information in Z is redundant and the pairing will be ineffective. This corresponds to the measure R_a^2 of Section 2.2 in the following way. We expect that the correlation between Y and X to be greater than Y and Z , (included covariates more effective than the omitted ones). In such a case, for the normal, linear model, R_a^2 is a monotonically decreasing function of ρ .

Each factor was inspected at two levels. In particular, x_3 was chosen to be the standardized versions of the following two vectors:

$(1, 1, \dots, 1, -1, -1, \dots, -1, 1, 1, \dots, 1, -1, -1, \dots, -1)$, and

$$(1, 2, \dots, \frac{m}{2}, 1, 2, \dots, \frac{m}{2}, 1, 2, \dots, \frac{m}{2}, 1, 2, \dots, \frac{m}{2})$$

(Note that z was also standardized in the computer trial.)

Based upon the asymptotic calculations, $MSE(\hat{\tau}_{\logit})$ was uniformly smaller than $MSE(\hat{\tau}_{ML})$, which suggests that logistic regression is the preferred analysis for large sample sizes. Also, of the seven factors considered, it was found that β_1 and x_3 had a much smaller influence on the mean square errors of $\hat{\tau}_{ML}$ and $\hat{\tau}_{\logit}$.

The next step was to determine the extent to which the asymptotic results apply for finite sample sizes. Retaining only the five factors that were deemed important in the calculations, a 2^3 by 3^2 simulation experiment was carried out and mean square errors were simulated using GAUSS (Edlefsen and Jones, 1986). (See Appendix 1). For each of the 72 runs, 3000 replications were performed, and the different factor levels chosen were:

$$\begin{aligned} m &= 30, 60, 100, & \beta_2 &= 0.4, 1.2, & \rho &= 0.2, 0.8, \\ \gamma &= 0, 0.5, 1.2, & \beta_3 &= 0.3, 1.2, \end{aligned}$$

(x_3 was set to be the standardized form of $(1, 1, \dots, 1, -1, -1, \dots, -1, 1, 1, \dots, 1, -1, -1, \dots, -1)$, and β_1 was set to zero.)

As $\hat{\tau}_{ML}$ can give rise to serious bias for small sample sizes, an additional matched pairs estimator was included in the comparison with $\hat{\tau}_{\logit}$. This was $\hat{\tau}_H = \ln\left(\frac{n_{10} + 0.5}{n_{01} + 0.5}\right)$, which seemed to perform better in Jewell's (1984) comparison of several sample estimators of the log odds ratio.

Tables 2 and 3 indicate the relative MSE performance of $\hat{\tau}_{ML}$ and $\hat{\tau}_H$ to $\hat{\tau}_{logit}$ by listing the values of r_{ML} and r_H for each of the 72 runs, where

r_{ML} is the median of 3000 replicated values of $\frac{(\hat{\tau}_{logit} - 2\beta_2)^2}{(\hat{\tau}_{ML} - 2\beta_2)^2}$, and

r_H is the median of 3000 replicated values of $\frac{(\hat{\tau}_{logit} - 2\beta_2)^2}{(\hat{\tau}_H - 2\beta_2)^2}$.

Medians were used as a robust measure since forming ratios sometimes produced extreme values.

1) Most of the values of r_{ML} and r_H are less than 1, which indicate that logistic regression performs better, with the exception when ρ is small and γ is large. This agrees with intuition since for our regression model, the smaller ρ is the less "information" about the omitted covariate is being captured, and the larger γ is the more bias we expect in estimating the treatment effect.

2) The most reliable factor determining the behavior of r_{ML} and r_H turns out to be the sample size. As the number of pairs increases, the values of r_{ML} and r_H will eventually drop below 1. This holds even for the anomaly observed for the small ρ large γ where further simulation showed that more than 200 pairs were needed.

3) In general as gamma increases, r_{ML} and r_H also increase.

4) As β_2 increases both r_{ML} and r_H decrease, while as β_3 increases both r_{ML} and r_H increase, although these two trends seemed reversed when ρ is small and γ is large.

5) r_{ML} and r_H do not seem to decrease as ρ increases unless γ is large.

We also note that overall there are no sharp differences between the behavior of r_{ML} and $r_{H'}$, aside from the observation that, for small sample sizes, increasing β_3 causes r_H to decrease when β_2 is large, and increasing γ causes r_H to decrease when β_2 and ρ are large.

4. CONCLUSION

The decision to break the matches for analysis of a matched pairs design and model the pairing will depend on whether the regression model is normal linear or Bernoulli logistic. With the normal linear model, the estimators of the treatment effect from the regression analysis and the matched pairs analysis are both unbiased but they have different variances. Regression turns out to be a poor alternative to matched pairs analysis (or the paired t-test) unless the number of pairs is very small and the measured covariates are able to account for a large proportion of the "information" of the omitted covariates. When the model is Bernoulli logistic, analytical results are only available asymptotically, and these indicate that logistic regression will in general be more efficient than the matched pairs analysis. For finite sample sizes, simulations suggest that logistic regression is still more favorable, but that the matched pairs analysis is preferable when the number of pairs is small, and when the regression coefficient of the omitted covariate is large and the correlation between the known covariate and the omitted covariate is small.

Table 2) Values of r_{ML}^*

		$\rho = 0.2$					
		$\gamma = 0$		0.5		1.2	
		β_3		β_3		β_3	
m	β_2	0.3	1.2	0.3	1.2	0.3	1.2
30	0.4	1.00	1.00	0.99	0.99	0.97	0.96
	1.2	0.76	0.99	0.77	1.07	1.52	1.65
60	0.4	0.95	0.97	0.96	0.96	1.01	1.00
	1.2	0.59	0.69	0.66	0.81	1.80	1.33
100	0.4	0.95	0.96	0.96	0.99	1.35	1.09
	1.2	0.53	0.66	0.62	0.72	2.59	1.51

		$\rho = 0.8$					
		$\gamma = 0$		0.5		1.2	
		β_3		β_3		β_3	
m	β_2	0.3	1.2	0.3	1.2	0.3	1.2
30	0.4	1.00	1.00	1.00	1.00	0.97	1.00
	1.2	0.79	1.01	0.83	1.04	1.07	0.99
60	0.4	0.93	0.96	0.96	0.98	0.97	1.00
	1.2	0.60	0.69	0.64	0.81	0.79	1.03
100	0.4	0.95	0.97	0.96	0.99	0.96	0.98
	1.2	0.54	0.71	0.58	0.74	0.76	0.84

KEY: $m \equiv$ number of pairs, $\beta_2 \equiv$ regression coefficient for considering the treatment effect, $\beta_3 \equiv$ regression coefficient for the known covariate, $\gamma \equiv$ regression coefficient for the omitted covariate, and $\rho \equiv$ correlation between x_3 and z (the omitted covariate vector).

* r_{ML} is defined as the median of 3000 replicated values of $\frac{(\hat{\tau}_{logit} - 2\beta_2)^2}{(\hat{\tau}_{ML} - 2\beta_2)^2}$.

Table 3) Values of r_H^*

		$\rho = 0.2$					
		$\gamma = 0$		0.5		1.2	
		β_3		β_3		β_3	
m	β_2	0.3	1.2	0.3	1.2	0.3	1.2
30	0.4	1.00	1.00	1.00	1.00	1.00	1.00
	1.2	1.03	0.75	1.08	0.75	1.72	1.13
60	0.4	0.97	1.00	0.97	0.99	1.03	1.00
	1.2	0.59	0.83	0.68	0.89	1.71	1.45
100	0.4	0.95	0.97	0.96	0.99	1.33	1.09
	1.2	0.59	0.72	0.64	0.75	2.47	1.56

		$\rho = 0.8$					
		$\gamma = 0$		0.5		1.2	
		β_3		β_3		β_3	
m	β_2	0.3	1.2	0.3	1.2	0.3	1.2
30	0.4	1.00	1.02	1.00	1.04	1.00	1.00
	1.2	1.14	0.74	0.98	0.60	0.84	0.58
60	0.4	0.95	0.97	0.96	1.00	1.00	1.00
	1.2	0.62	0.76	0.67	0.81	0.89	0.98
100	0.4	0.96	0.98	0.96	0.99	0.97	0.99
	1.2	0.57	0.74	0.63	0.81	0.79	0.91

$*r_H$ is defined as the median of 3000 replicated values of $\frac{(\hat{\tau}_{\log t} - 2\beta_2)^2}{(\hat{\tau}_H - 2\beta_2)^2}$.

ACKNOWLEDGEMENTS

We wish to thank Professor Thomas Santner for reviewing the manuscript. This research was supported in part by Hatch grant 151-406.

REFERENCES

- Breslow, N. E. and Day, N. E. (1980). *Statistical Methods in Cancer Research, Vol. 1: The Analysis of Case-Control Studies*. IARC Scientific Publications No. 32, International Agency for Research on Cancer, Lyon, France.
- Edlefsen, L. E. and Jones, S. D. (1986). *GAUSS*, documentation version 1.00, software version 1.49B. Washington: Aptech Systems Inc.
- Gail, M.; Wieand, S.; and Piantadosi, S. (1984). "Biased estimates of treatment effect in randomized experiments with non-linear regressions and omitted covariates," *Biometrika*, 71, pp. 431-444.
- Gail, M.; Tan, W.; and Piantadosi, S. (1988). "Test of no treatment effect in randomized clinical trials," *Biometrika*, 75, pp. 57-64.
- Greenland, S. (1986). "Partial and marginal matching in case-control studies", in *Modern Statistical Methods in Chronic Disease Epidemiology* ,(pp. 35-49) edited by Moolgavkar, S. and Prentice, R. John Wiley & Sons, Inc. , New York.
- Jewell, N. P. (1984). "Small-sample bias of point estimators of the odds ratio from matched sets," *Biometrics*, 40, pp. 421-435.
- Kupper, L. L.; Karon, J. M.; Kleinbaum, D. G.; Morgenstern, H.; and Lewis, D. K. (1981). "Matching in epidemiologic studies: validity and efficiency considerations," *Biometrics*, 37, pp. 271-292.
- McKinlay, S. M. (1977). "Pair matching -- a reappraisal of a popular technique," *Biometrics*, 33, pp. 725-735.

- Prentice, R. (1976). "Use of the logistic model in retrospective studies," *Biometrics*, 32, pp. 599-606.
- Royall, R. M. (1986). "Model robust confidence intervals using maximum likelihood estimators," *International Statistical Review*, 54, pp. 221-226.
- Rubin, D. B. (1973). "The use of matched sampling and regression adjustment to remove bias in observational studies," *Biometrics*, 29, pp. 185-203.

APPENDIX 1

The following tables list the simulated mean square errors of three estimators of the log odds ratio, with their respective standard errors included underneath in parentheses. $\hat{\tau}_{\text{logit}}$ is the estimator obtained from logistic regression, whereas $\hat{\tau}_{\text{ML}}$ and $\hat{\tau}_{\text{H}}$ are the two different matched pairs estimators. n1 is the number of replications in each run that led to undefined estimators for the matched pairs analysis, and n2 is the number of replications in each run that did not produce converging solutions for the logistic regression. (Note that the values for the case when m=60 have been omitted in the following tables to shorten the presentation.)

m=30								
$\hat{\tau}_{\text{ML}}$	$\hat{\tau}_{\text{H}}$	$\hat{\tau}_{\text{logit}}$	n1	n2	ρ	γ	β_2	β_3
.3656 (.0108)	.2858 (.0076)	.3265 (.0095)	18	0	0.2	0	0.4	0.3
.4857 (.0134)	.3625 (.0098)	.5039 (.0153)	82	14			0.4	1.2
.3581 (.0085)	.3666 (.0102)	.5507 (.0167)	587	7			1.2	0.3
.3710 (.0123)	.4707 (.0140)	.5565 (.0142)	880	216			1.2	1.2
.4030 (.0120)	.3111 (.0086)	.3135 (.0091)	30	0		0.5	0.4	0.3
.4694 (.0125)	.3496 (.0091)	.4755 (.0142)	90	15			0.4	1.2
.3269 (.0086)	.3508 (.0103)	.4783 (.0151)	662	8			1.2	0.3
.3870 (.0128)	.5081 (.0144)	.5049 (.0122)	893	210			1.2	1.2
.4657 (.0129)	.3544 (.0094)	.2819 (.0073)	71	0		1.2	0.4	0.3
.5151 (.0134)	.3762 (.0100)	.3852 (.0113)	132	2			0.4	1.2

m=30								
$\hat{\tau}_{ML}$	$\hat{\tau}_H$	$\hat{\tau}_{logit}$	n1	n2	ρ	γ	β_2	β_3
.3544 (.0111)	.4484 (.0129)	.5308 (.0117)	812	2	0.2	1.2	1.2	0.3
.4122 (.0141)	.5864 (.0156)	.5601 (.0123)	1028	84			1.2	1.2
.3804 (.0112)	.2979 (.0079)	.3296 (.0090)	20	0	0.8	0	0.4	0.3
.4698 (.0126)	.3491 (.0091)	.4877 (.0139)	74	13			0.4	1.2
.3353 (.0081)	.3341 (.0098)	.5397 (.0169)	596	6			1.2	0.3
.3757 (.0119)	.4739 (.0136)	.5444 (.0136)	798	195			1.2	1.2
.4071 (.0118)	.3185 (.0087)	.3544 (.0103)	20	0		0.5	0.4	0.3
.5371 (.0139)	.3931 (.0103)	.5970 (.0160)	174	62			0.4	1.2
.3292 (.0087)	.3667 (.0105)	.5291 (.0159)	647	38			1.2	0.3
.4567 (.0155)	.6418 (.0169)	.5054 (.0129)	1022	459			1.2	1.2
.4957 (.0127)	.3707 (.0092)	.4318 (.0120)	83	7		1.2	0.4	0.3
.6049 (.0169)	.4512 (.0132)	.6518 (.0176)	363	223			0.4	1.2
.3883 (.0127)	.5120 (.0143)	.5144 (.0121)	891	127			1.2	0.3
.6469 (.0228)	.9583 (.0235)	.5662 (.0174)	1374	808			1.2	1.2

m=100								
$\hat{\tau}_{ML}$	$\hat{\tau}_H$	$\hat{\tau}_{logit}$	n1	n2	ρ	γ	β_2	β_3
.1011 (.0029)	.0950 (.0026)	.0903 (.0023)	0	0	0.2	0	0.4	0.3
.1469 (.0045)	.1341 (.0039)	.1304 (.0035)	0	0			0.4	1.2
.2676 (.0091)	.2048 (.0058)	.1259 (.0036)	13	0			1.2	0.3
.3187 (.0095)	.2287 (.0057)	.1953 (.0064)	37	0			1.2	1.2

m=100								
$\hat{\tau}_{ML}$	$\hat{\tau}_H$	$\hat{\tau}_{logit}$	n1	n2	ρ	γ	β_2	β_3
.1123 (.0035)	.1047 (.0031)	.0877 (.0023)	0	0	0.2	0.5	0.4	0.3
.1443 (.0042)	.1316 (.0037)	.1211 (.0032)	0	0			0.4	1.2
.2700 (.0089)	.2054 (.0056)	.1232 (.0031)	17	0			1.2	0.3
.3268 (.0096)	.2369 (.0057)	.1972 (.0061)	61	0			1.2	1.2
.1368 (.0044)	.1256 (.0038)	.1014 (.0024)	0	0		1.2	0.4	0.3
.1679 (.0053)	.1508 (.0044)	.1156 (.0029)	0	0			0.4	1.2
.3281 (.0098)	.2340 (.0058)	.3439 (.0056)	45	0			1.2	0.3
.3527 (.0094)	.2471 (.0057)	.2837 (.0059)	89	0			1.2	1.2
.0979 (.0029)	.0925 (.0026)	.0882 (.0023)	0	0	0.8	0	0.4	0.3
.1478 (.0050)	.1345 (.0042)	.1315 (.0037)	0	0			0.4	1.2
.2604 (.0087)	.1982 (.0055)	.1220 (.0036)	13	0			1.2	0.3
.3231 (.0097)	.2298 (.0057))	.2032 (.0062)	38	0			1.2	1.2
.1089 (.0039)	.1011 (.0033)	.0925 (.0025)	0	0		0.5	0.4	0.3
.1834 (.0056)	.1631 (.0046)	.1641 (.0044)	0	0			0.4	1.2
.2858 (.0093)	.2115 (.0058)	.1418 (.0044)	22	0			1.2	0.3
.3565 (.0094)	.2534 (.0056)	.2753 (.0090)	101	14			1.2	1.2
.1533 (.0048)	.1389 (.0041)	.1154 (.0031)	0	0		1.2	0.4	0.3
.2789 (.0089)	.2322 (.0067)	.2266 (.0071)	0	0			0.4	1.2
.3304 (.0097)	.2333 (.0056)	.1805 (.0049)	54	1			1.2	0.3
.3624 (.0078)	.2697 (.0059)	.3552 (.0096)	223	45			1.2	1.2

APPENDIX 2: SIMULATION DETAILS

The GAUSS program was run on a IBM personal system/2TM Model 60 machine. 3000 replications were performed for each combination of the five parameter values: m , ρ , γ , β_2 , and β_3 . When calculating $\hat{\tau}_{logit}$, replications that did not give converging solutions within 30 iterations of a Newton-Raphson algorithm were discarded. (To ensure that 30 is a reasonable choice, a subset of the data sets that did not converge within 30 iterations were inspected and these were all verified to yield nonconvergent solutions.) For the matched pairs analysis, a replication was discarded when $\hat{\tau}_{ML}$ was undefined.